

RHIC Computing Facility

RHIC S&T Review

***Eric Lançon
August 23, 2016***



Outline

- Status of RCF, synergies with ATLAS Tier-1
- Performance in recent RHIC runs
- Future technological and data challenges
- Synergies with BNL Computing Initiative
- B725 infrastructure project

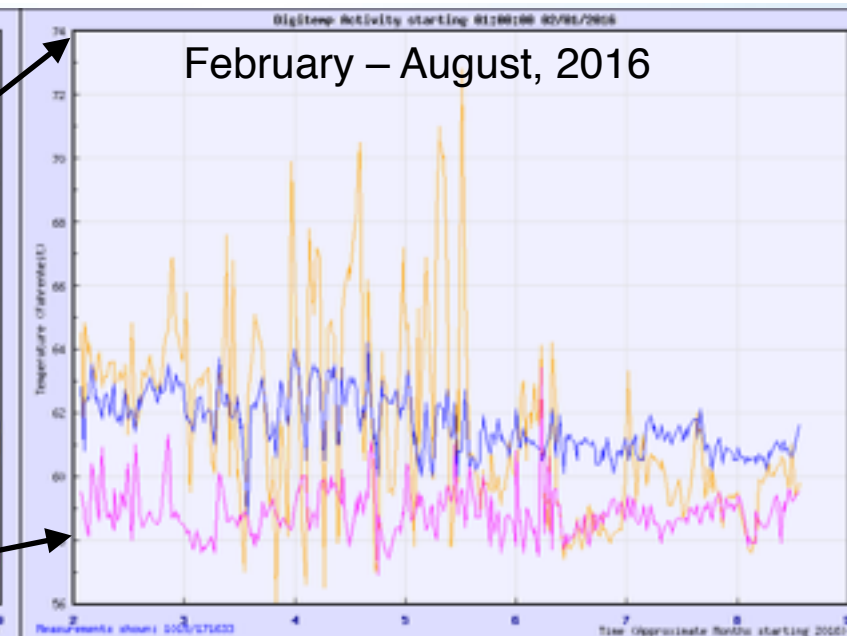
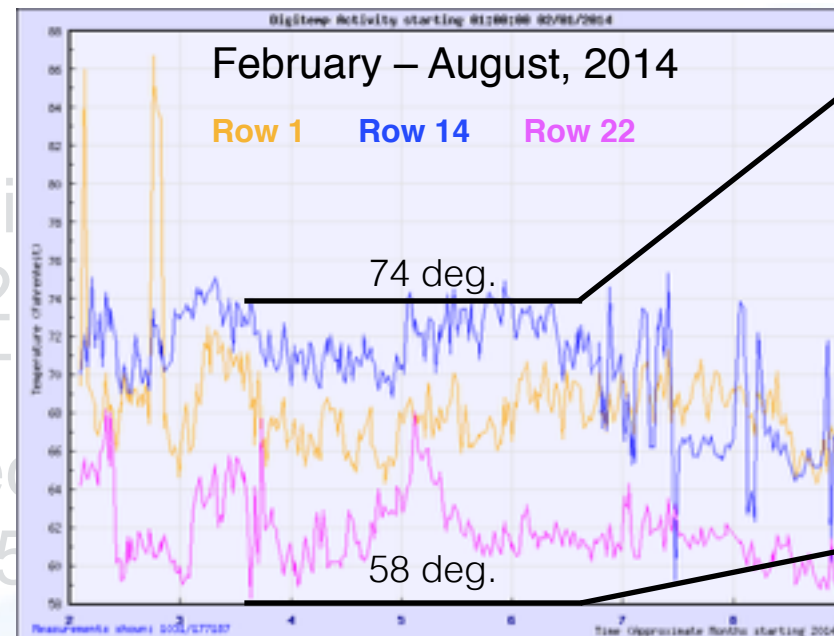
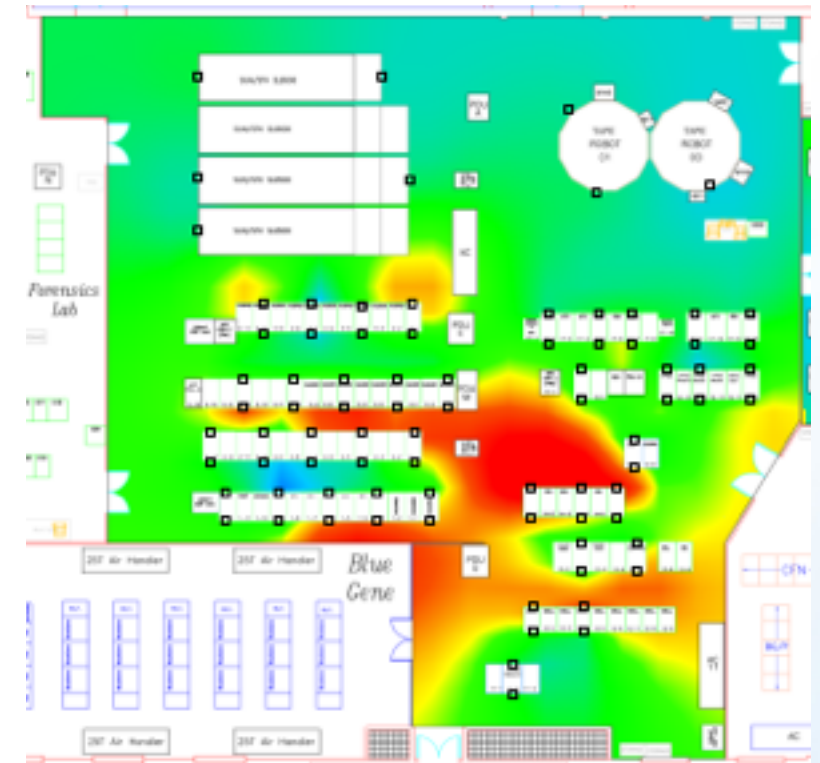
Recommendations from 2014 S&T review

1. BNL is encouraged to resolve the HVAC problem at RACF as soon as possible.

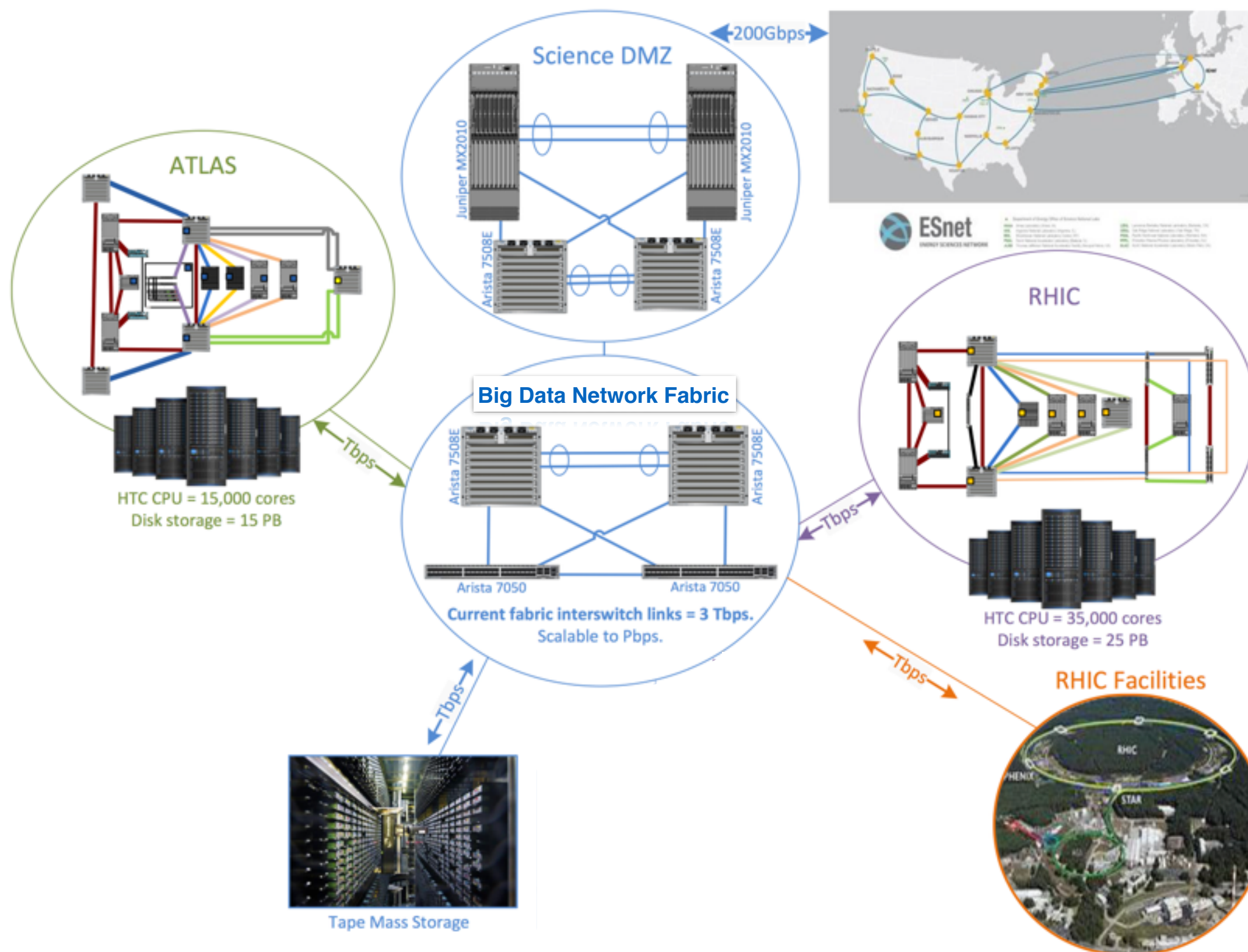
2. The RACF and the detector collaborations should analyze the processing capacity required to perform the necessary production runs (in units of HS06*year) compare to the available capacity of RACF. If additional capacity is required, a plan to acquire necessary capacity by 2015 should be developed. This plan should be submitted to DOE by January 1, 2015.

3 additional Heating & Ventilating Air Conditioning units installed reducing temperature level & gradient in computing room

Temperature map in computing room

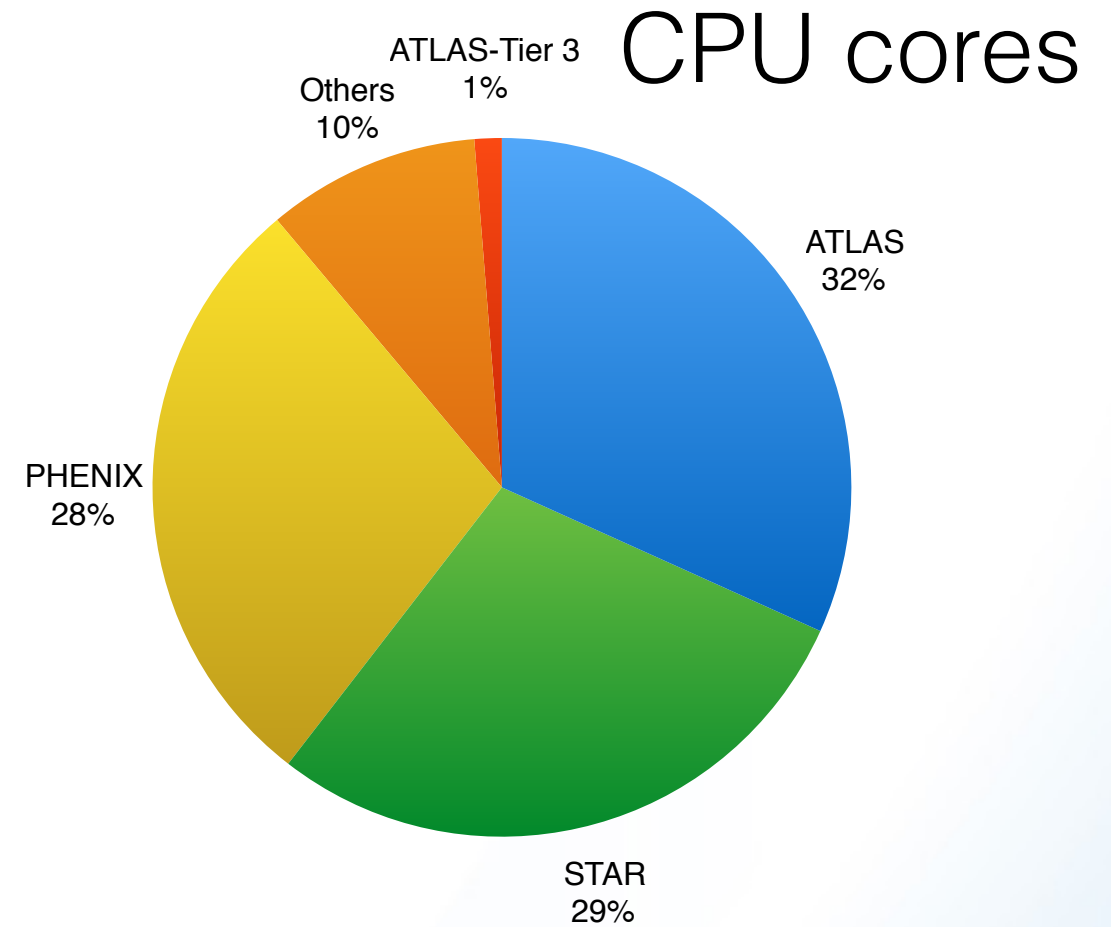


RCF today



Capacities as of today

- **55k CPU cores**
 - 3% HPC of capacity, will increase in the next months
- **~45 PB of disk storage**
 - of various technologies
- **~80 PB of tape storage**
 - 4th HPSS (High Performance Storage System) site worldwide
 - first within the US⁽¹⁾



Site	HPSS sites
	(ECMWF) European Centre for Medium-Range Weather Forecasts
	(NOAA-RD) National Oceanic and Atmospheric Administration Research & Development
	(UKMO) United Kingdom Met Office
	(BNL) Brookhaven National Laboratory
	(LBNL-User) Lawrence Berkley National Laboratory - User
	(LANL-Secure) Los Alamos National Laboratory - Secure
	(ORNL) Oak Ridge National Laboratory
	(NCAR) National Center for Atmospheric Research

(1) http://www.hpss-collaboration.org/learn_who_petabyte_data.shtml

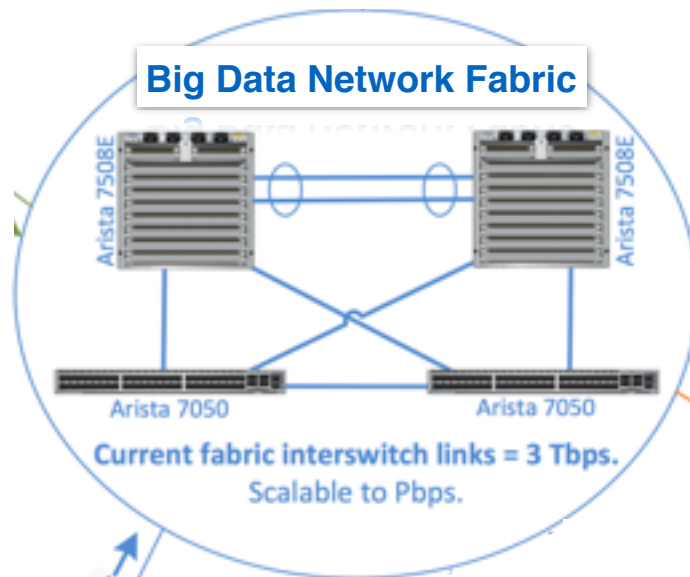
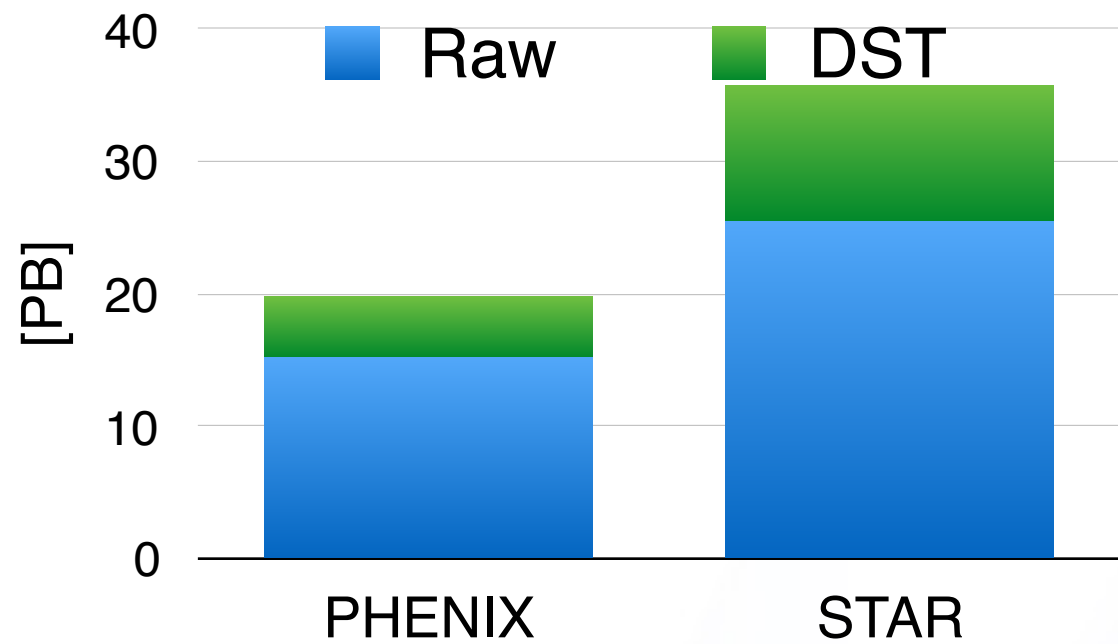
Status of RCF

- RCF performed well during 2016 run
- Resources are ~fully utilised
- Hardware (CPU) is getting old, migration to new tape generation needed (space in HPSS)
- Increase of resources needed in the next years

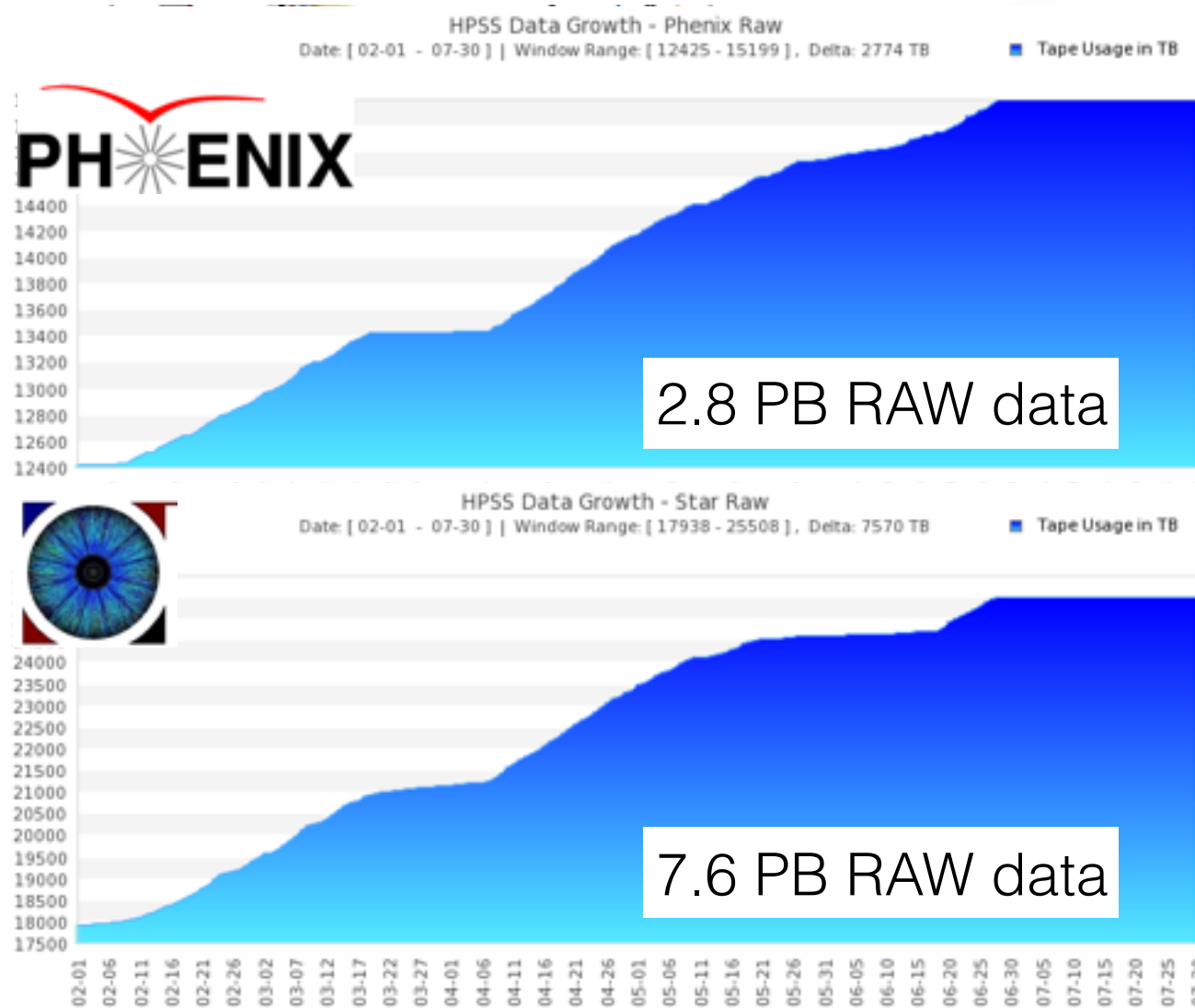
Performance in 2016

- No issue in writing RAW data to tape

HPSS Occupancy [PB]

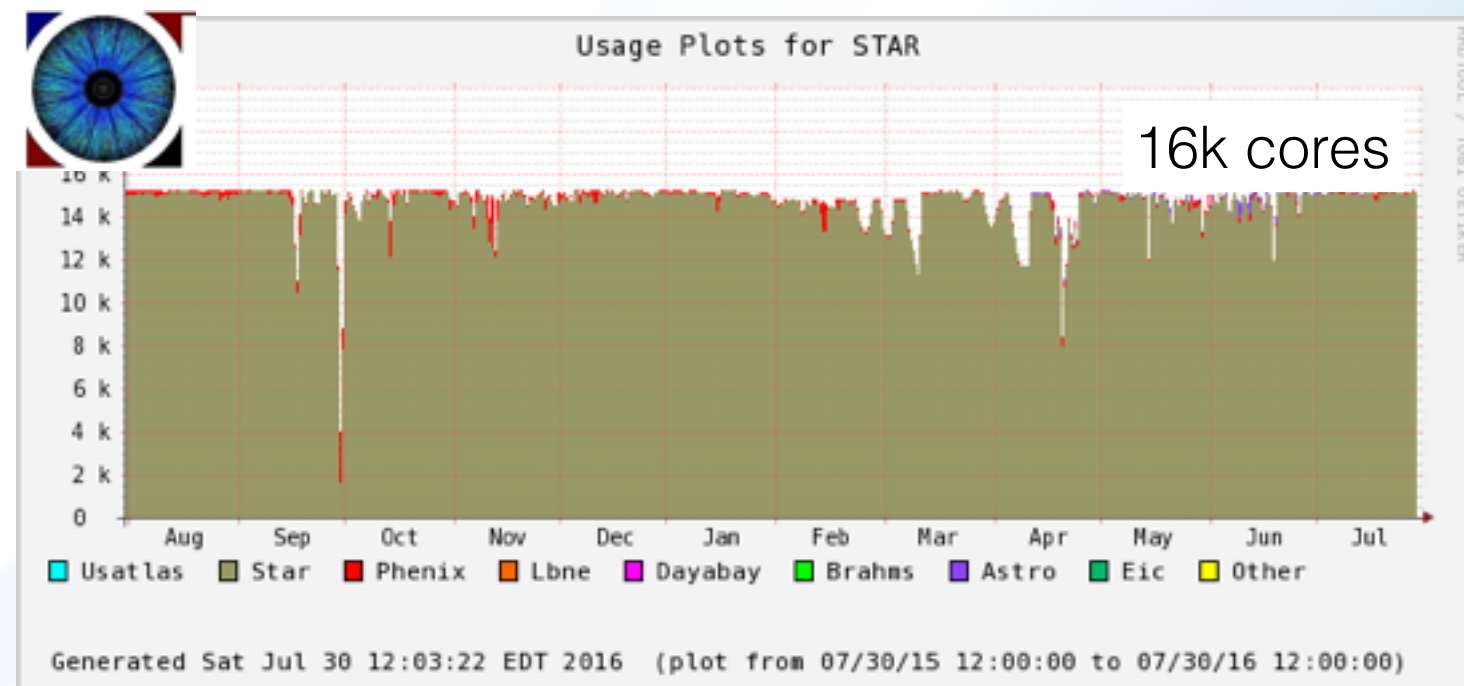
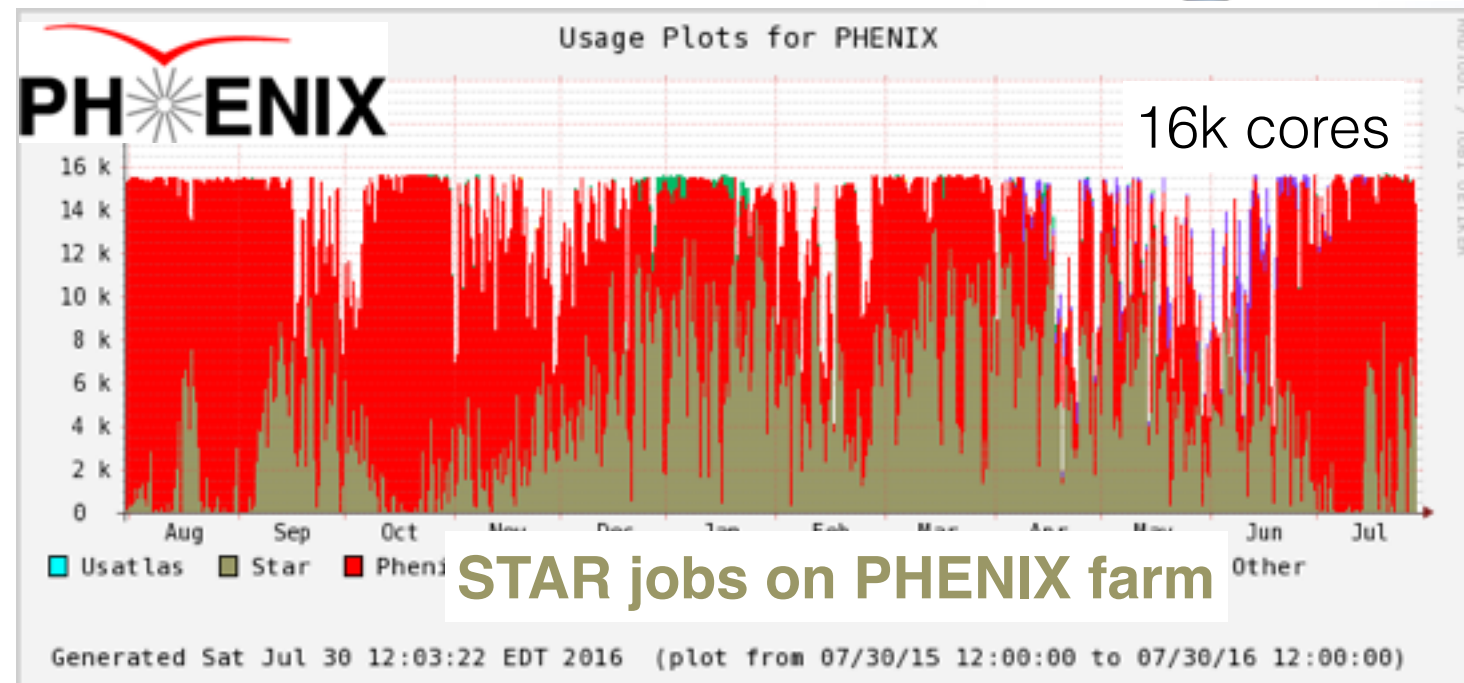
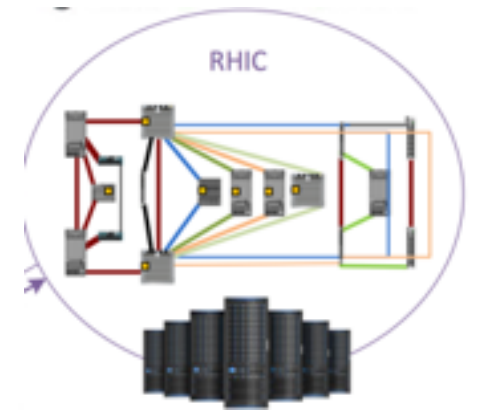


Tape Mass Storage



CPU usage of the farms

- Two distinct computing farms of equal size, one for PHENIX, one for STAR
- Storage distributed on computing nodes
 - Reconstruction jobs of experiment A cannot run on farm of experiment B
- STAR farm almost continuously saturated while PHENIX farm is not
- PHENIX farm used by STAR analysis jobs when no PHENIX activity
 - Optimisation of batch system (Condor) performed by RCF,
 - STAR analysis workflow optimisation to be done (too long jobs)
- *Lesson for the future*
 - *Computing models (workflow management, data organisation,...) and technological choices (storage, CPU,...) of experiments should not be too different in order to benefit from a global pool of resources*

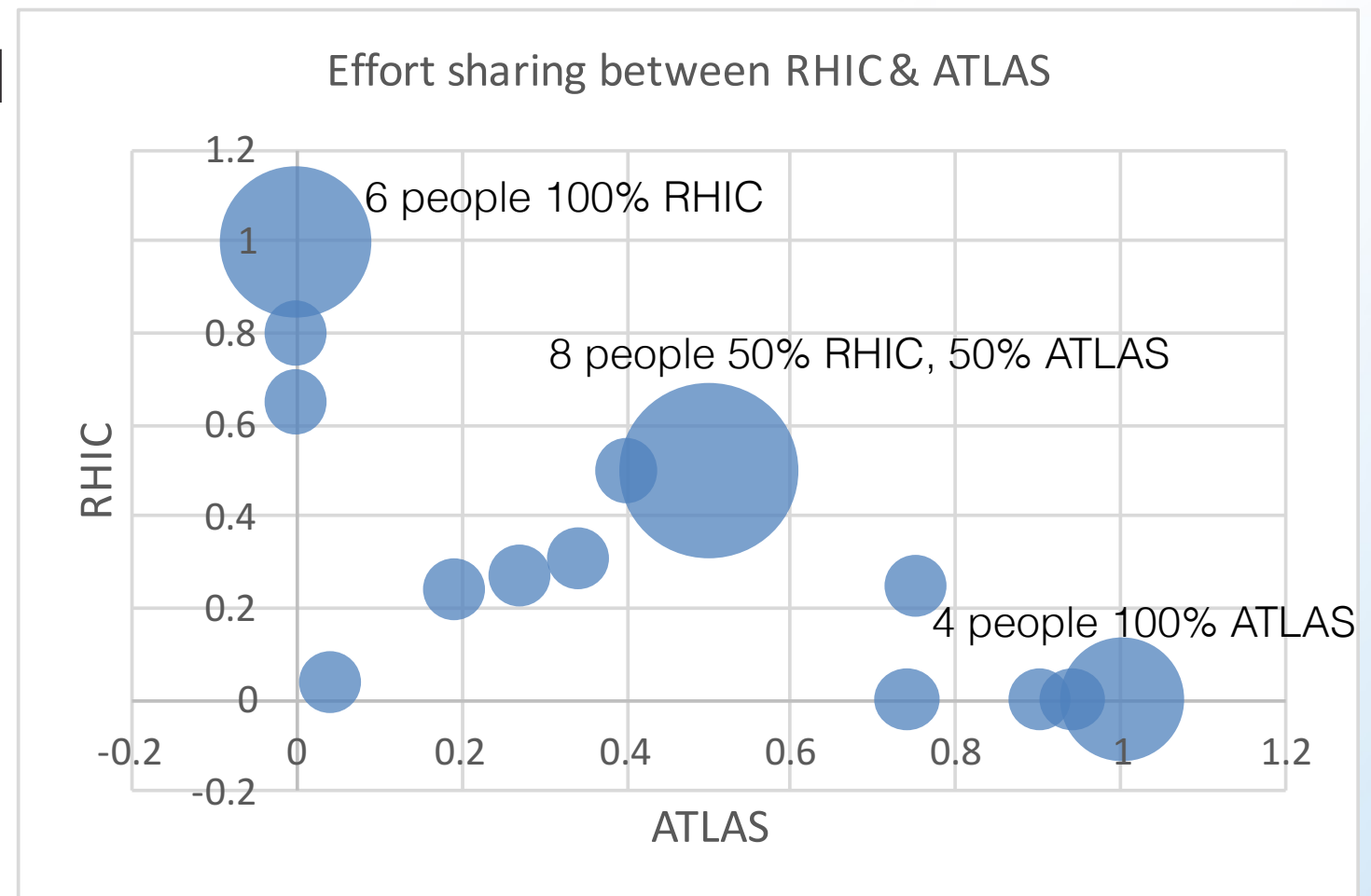


Synergy with ATLAS Tier-1

- Economy of scale (operation, purchase,...)
- Common procedures and configurations (resilience)
- Common tools (batch system, storage, network)
- Expertise from RCF benefits to ATLAS (and vis versa)
- Access to Grid and cloud computing expertise developed in ATLAS
- ...

Synergy with ATLAS Tier-1

- 13.6 FTE for RHIC
 - Support from ITD included
 - 6 people are 100% RHIC (storage, infrastructure, user support,...)
 - 8 people 50/50 (batch, system administration & configuration,...)
- About the right size of effort provided new RHIC experiments do not develop complex computing models

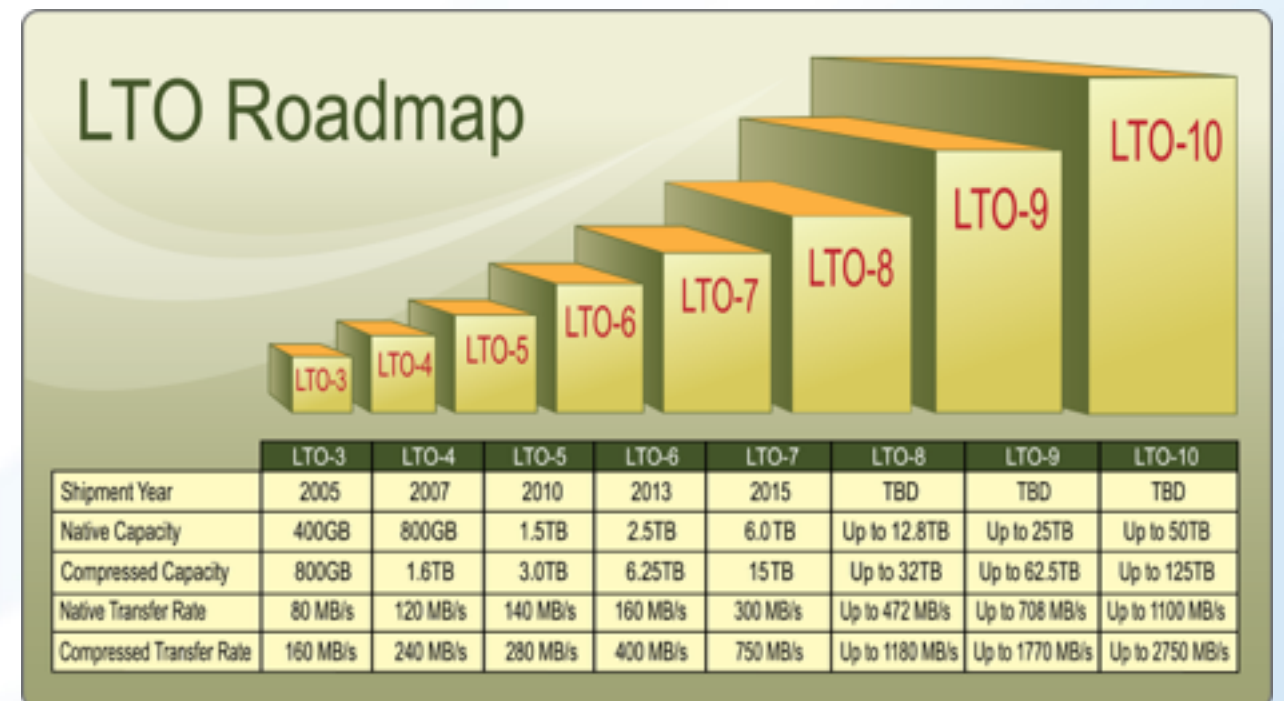
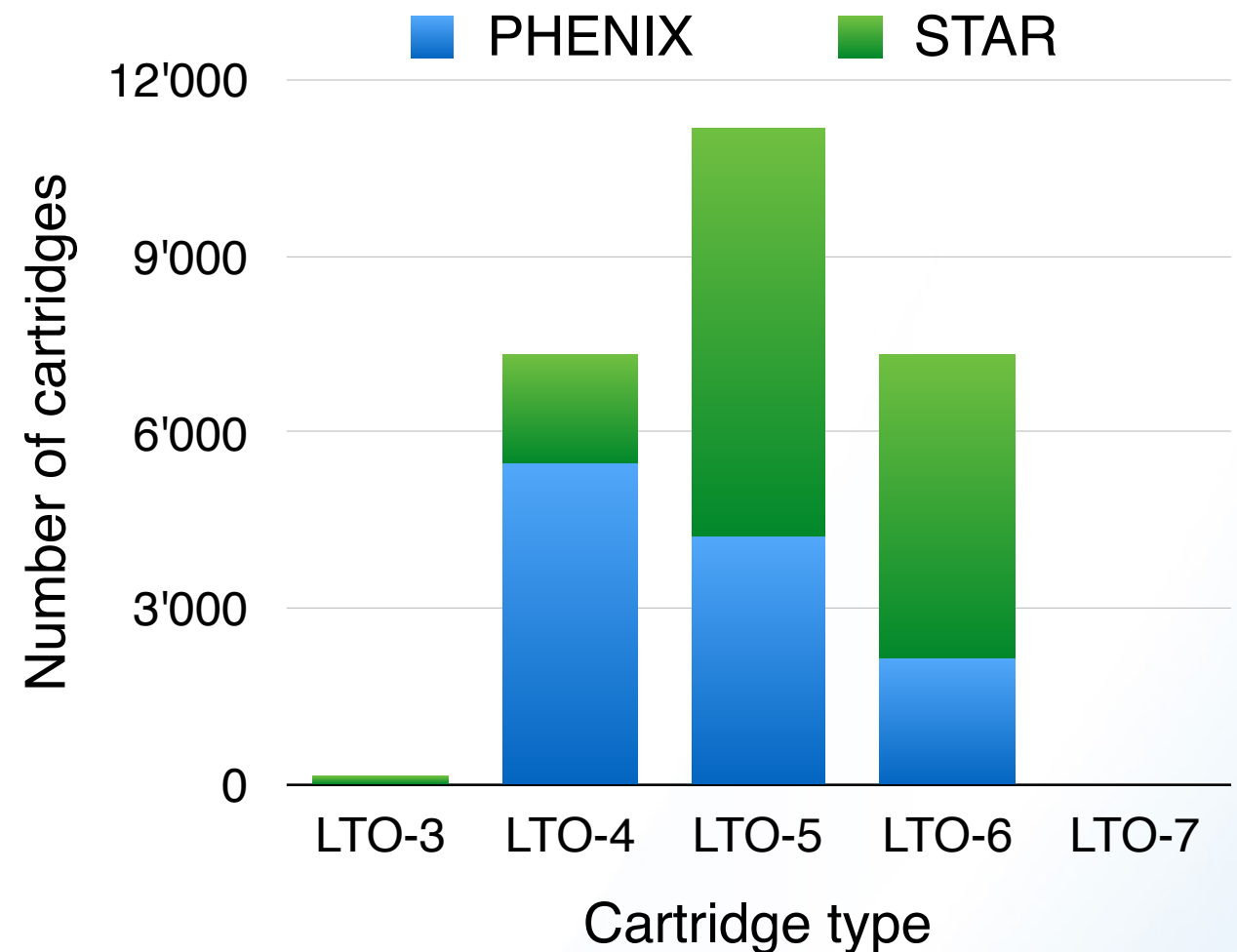


Future technological and data challenges

- **Future of computing is multi-core**
 - New hardware are multi-core 16, 32, 64,.... with less and less memory per core
 - Could software of RHIC future experiments be multi-core?
 - Is it worth the effort for existing experiments?
- **Object store technology**
 - ATLAS will migrate to Ceph (2-5 years)
 - To be considered for sPHENIX and eRHIC?
- **RHIC hardware is getting old, ~25% older than 5 years**
- **Tape technology**
 - 2 generations behind in tape technology
 - Only one copy of RAW data on tape
- **Data preservation**
 - Access to data and software maintenance long after data taking

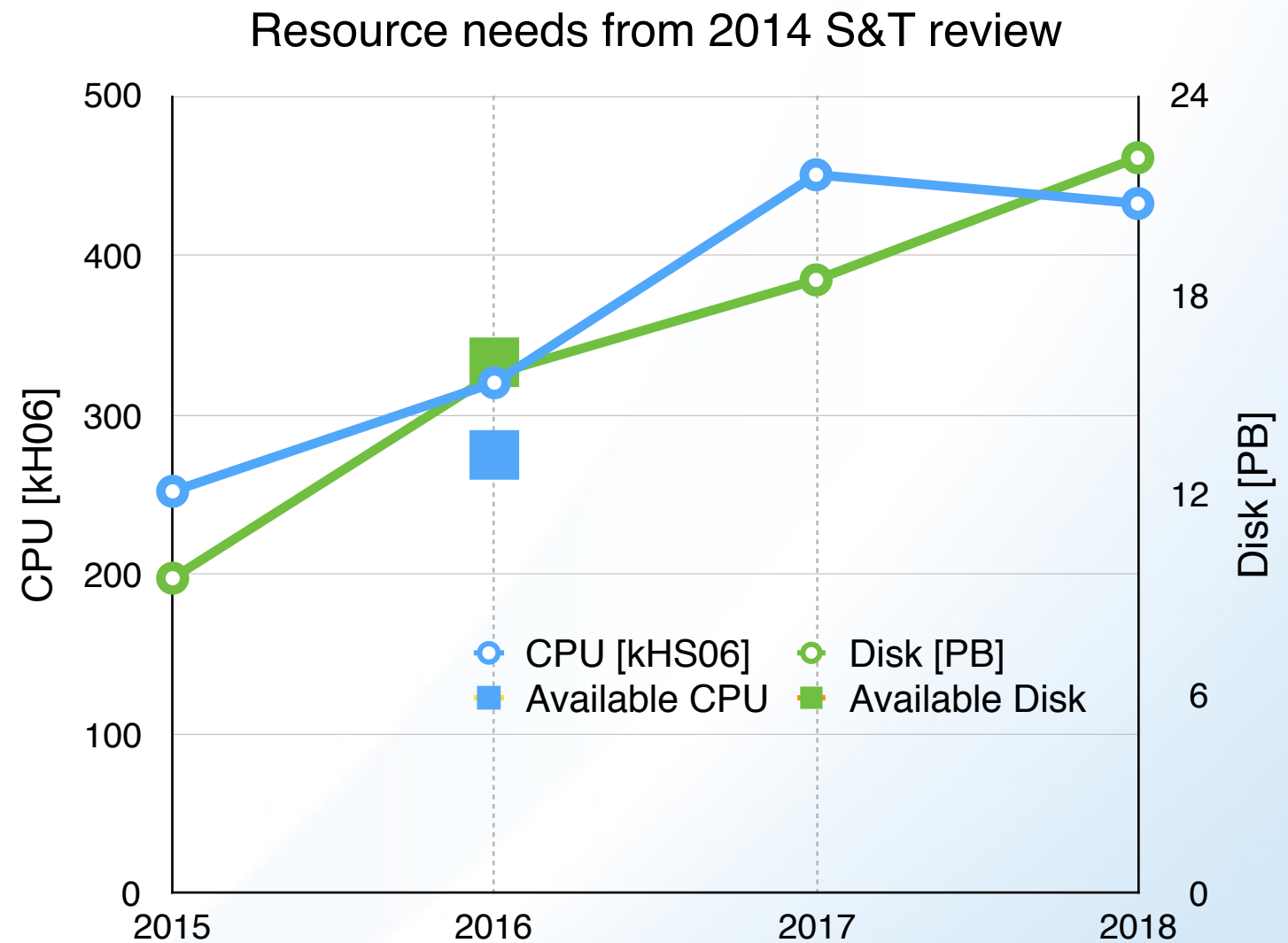
Tape migration

- Need to migrate archived data to new tape technology (LTO-7)
 - ~7 more capacity / tape
 - ~3 time faster
- LTO-7 tape drives cannot read LTO-4 and older types
 - Data on LTO-4 copied onto LTO-7
- 2 copies of RAW data will be made in the migration process
 - Today 1 copy of RAW data



CPU & Disk resources for next years

- Today's resources just match anticipated needs from 2014 S&T review
- 25% of capacity is older than 5 years and need to be replaced
- Projected 2017 needs (including replacement) ~1.9 current capacity
- Projection did not include running in 2017
- **Real 2017 needs ~2.1 current capacity**

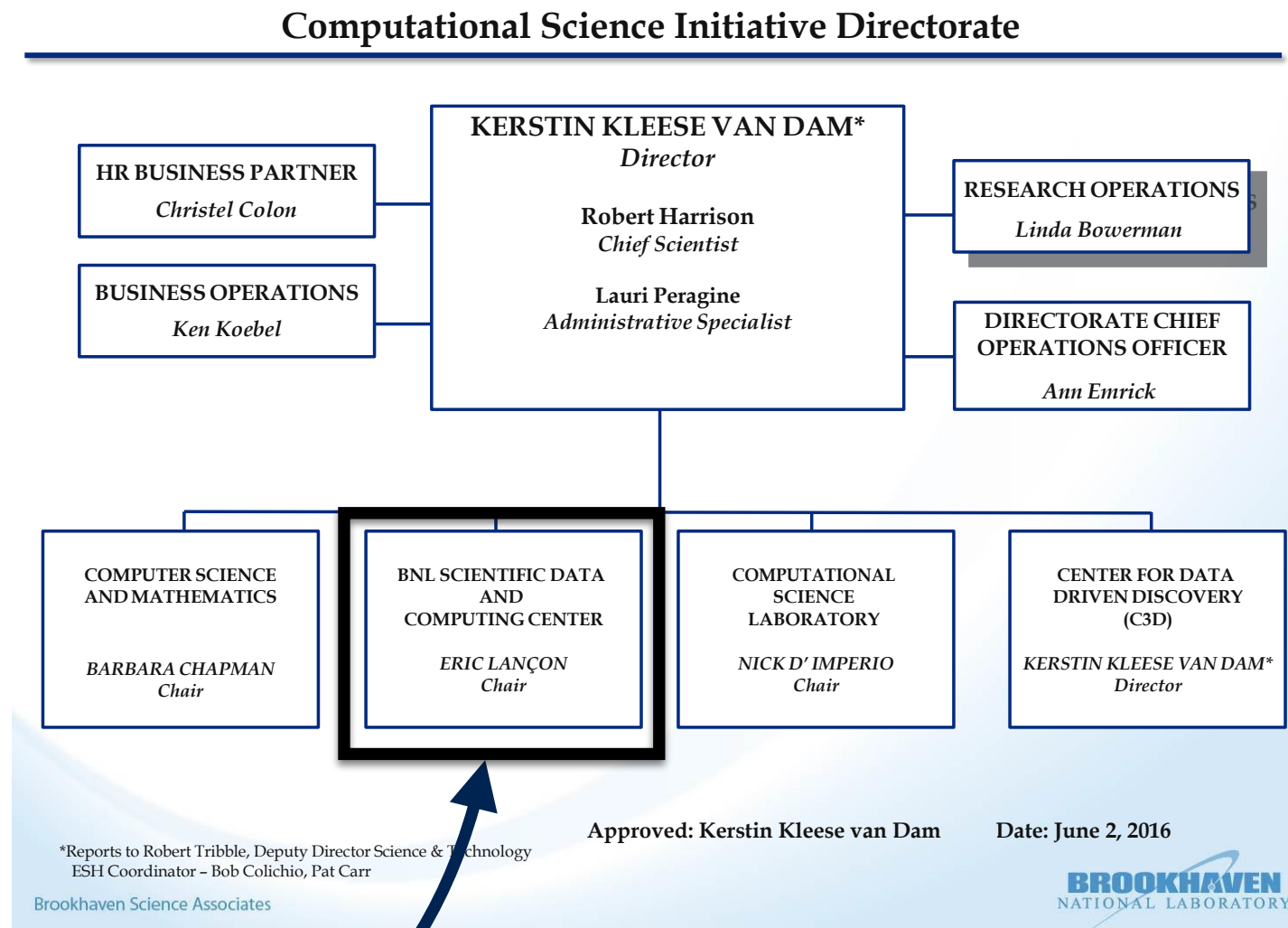


Computational Science Initiative : CSI

- **CSI**

- Leverage laboratory investments in scientific computing across multiple programs
- Patterns : universities (Columbia, Cornell, New York University, Stony Brook, and Yale) and companies including IBM Research.

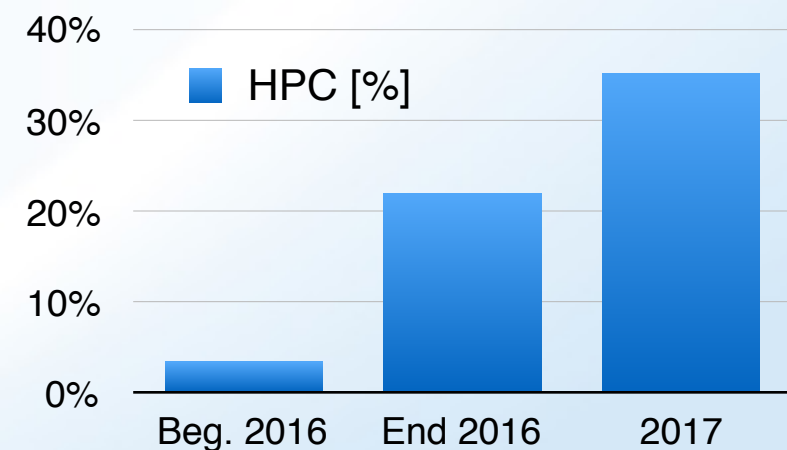
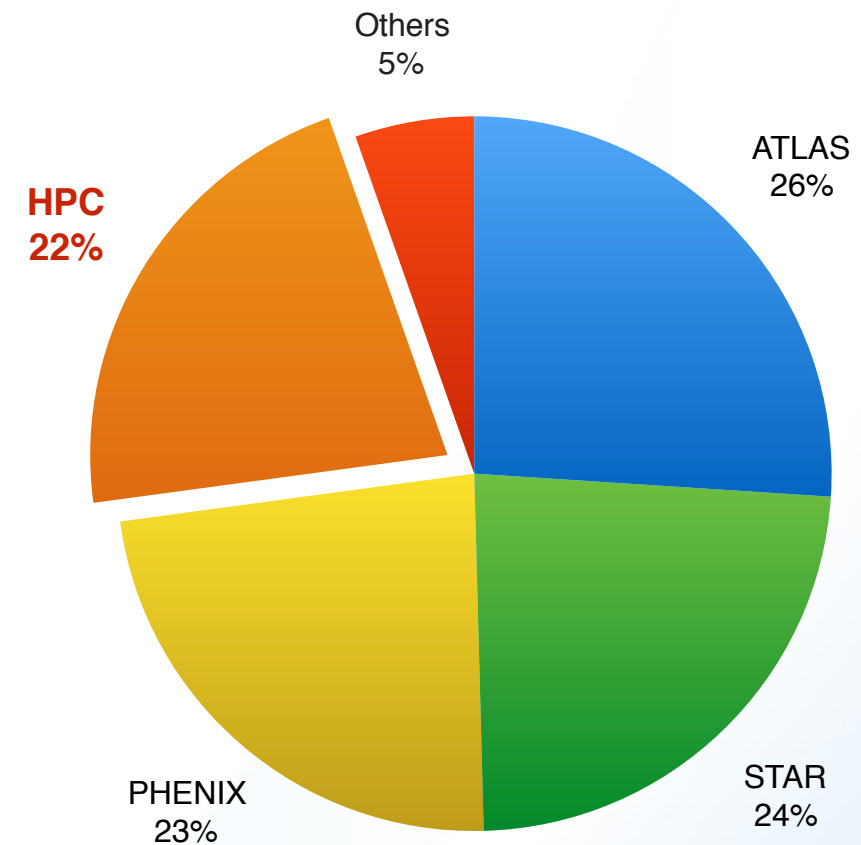
- **SDCC:** Scientific Data and Computing Center of CSI



RHIC and ATLAS Computing Facility operates SDCC

- **SDCC** is operated by **RACF**
- It includes components from
 - Laboratory's Institutional Cluster
 - CFN (Center for Functional Nano-materials)
 - Atmospheric Radiation Measurement
 - USQCD
 - ...
- Heavy investments by CSI in HPC

SDCC en 2016
70k cores

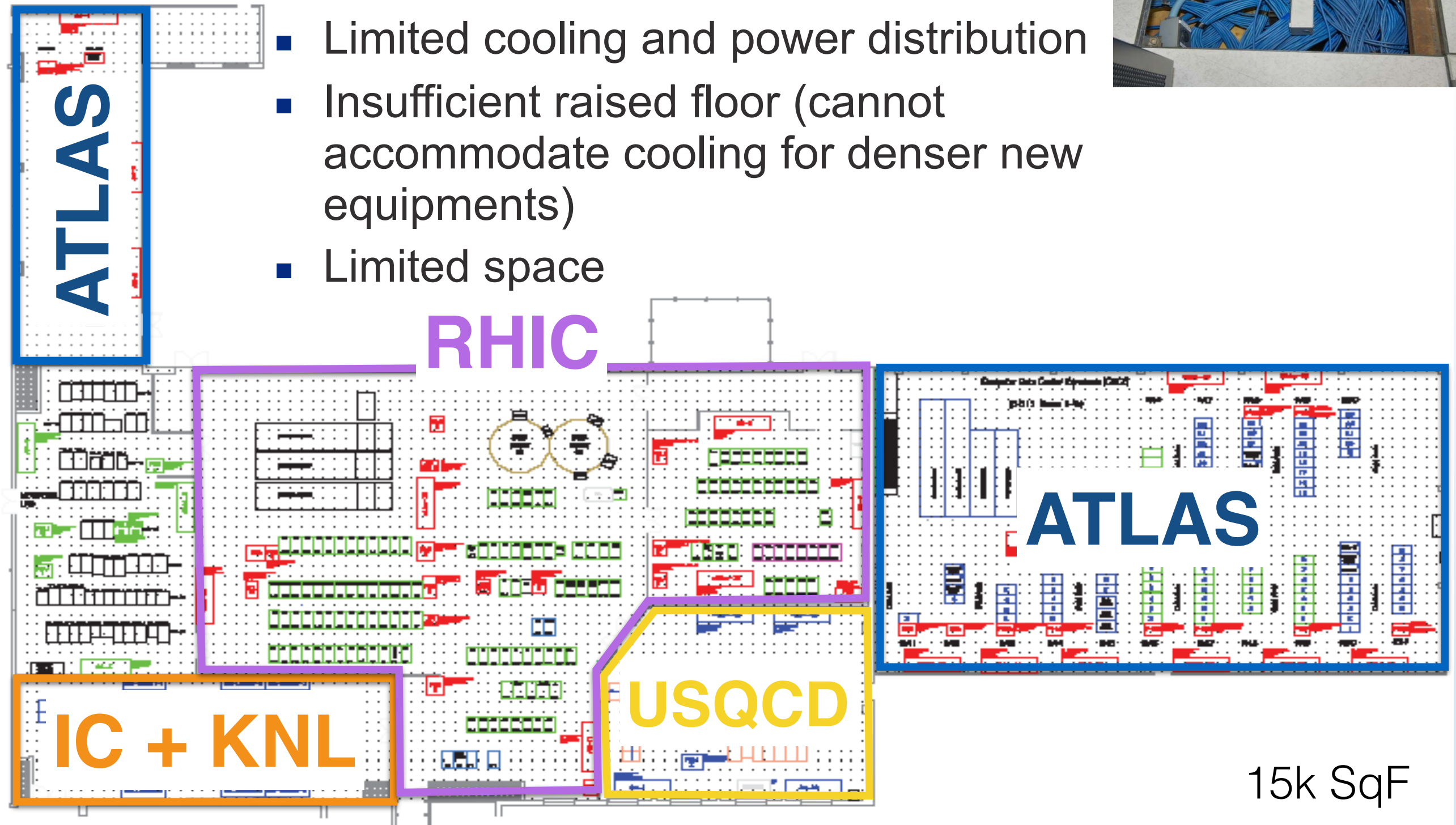


Synergies with BNL Computing Initiative

- CSI is purchasing or complementing purchases in the area of HPC computing (multi-core interconnected nodes)
 - Institutional cluster (Fall 2016, 2x 2017)
 - Knight Landings (KNL) Intel farm (Fall 2016). Initiated by BNL QCD group and RIKEN, CSI doubled the capacity
- These resources will be made available to RHIC program in opportunistic mode
 - May add 10% to RHIC resources?
 - Issue : manpower to port RHIC codes on KNL?
- Leverage on expertise in data processing, networking & storage technologies developed for RHIC and ATLAS

Computing room(s)

- OLD installations
- Limited cooling and power distribution
- Insufficient raised floor (cannot accommodate cooling for denser new equipments)
- Limited space

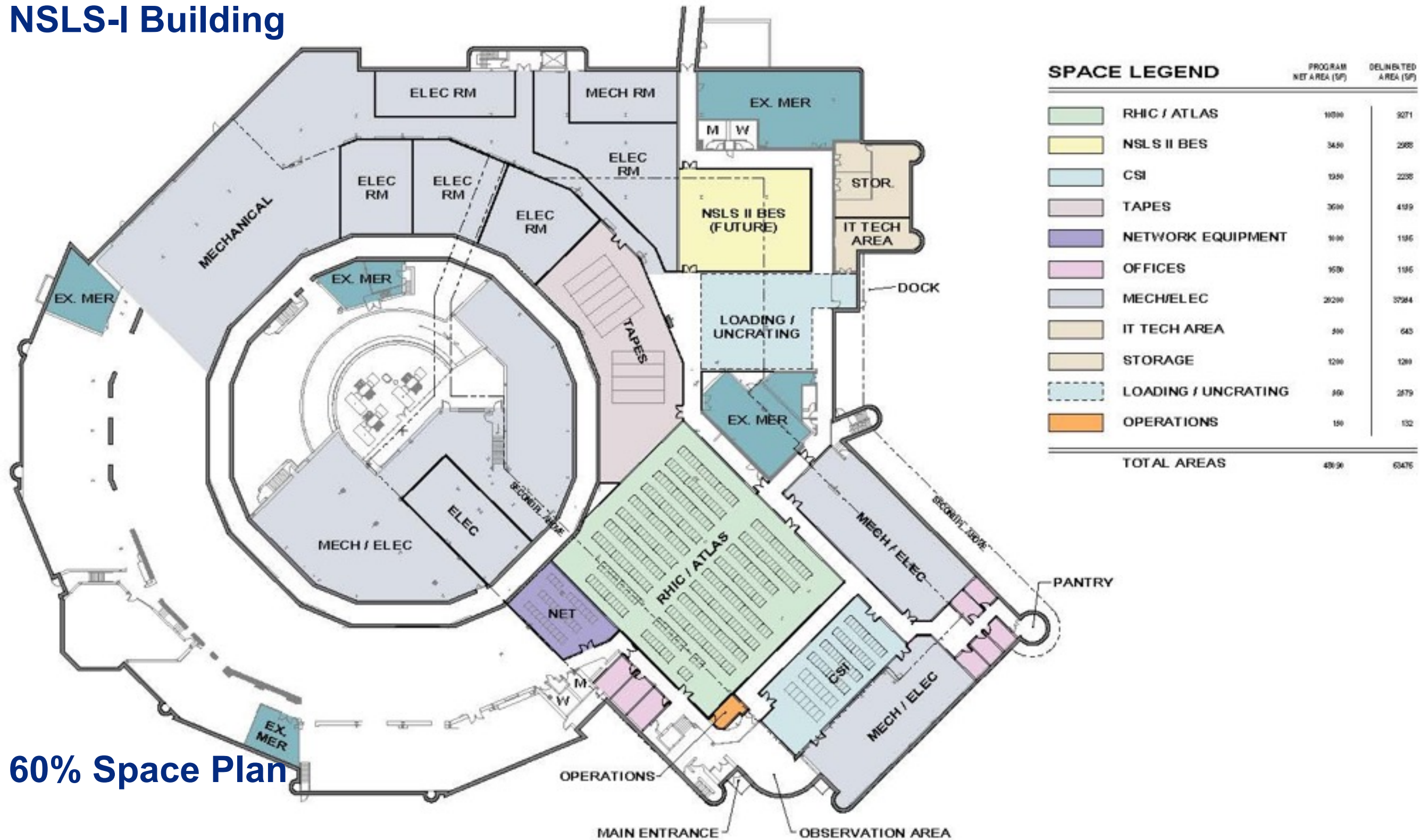


15k SqF

New computing room needed

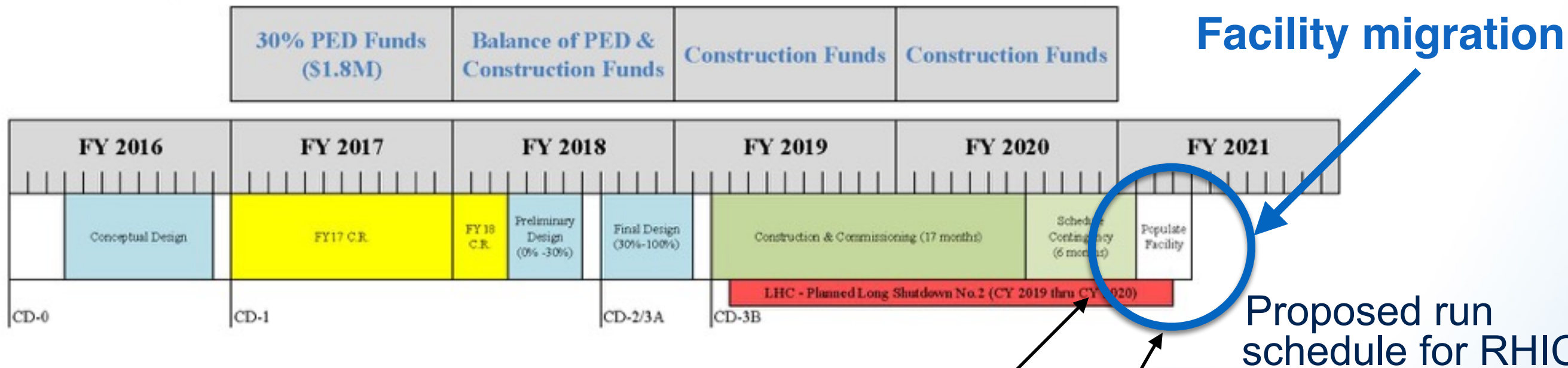
Core Facility Revitalisation – Conceptual Design

NSLS-I Building



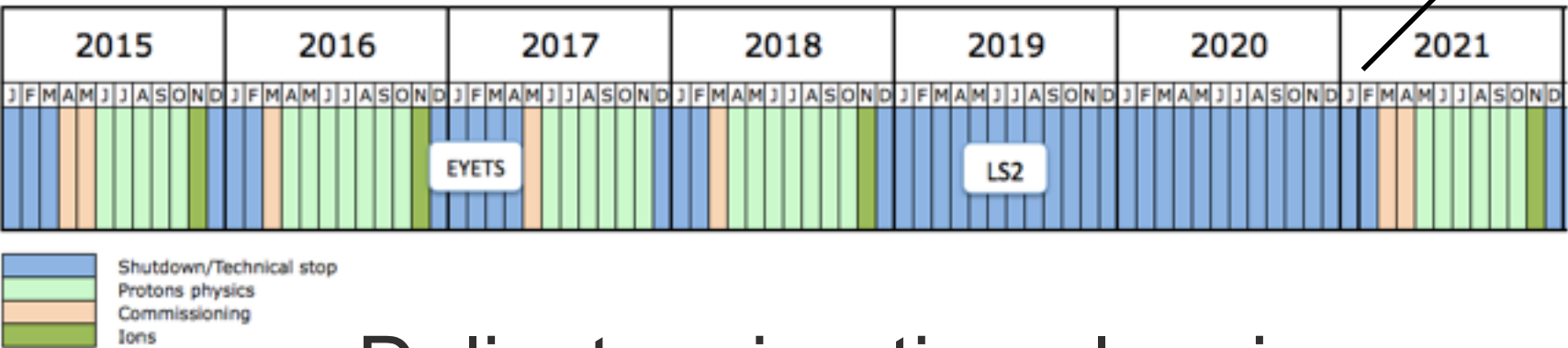
CFR – Preliminary Schedule

Preliminary CFR Funding Analysis - 1 Yr CR 2017 (Renovation Alternative)



Longer term LHC schedule

The outline LHC schedule out to 2035 presented by Frederick Bordry to the SPC and FC June 2015 can be found [here](#)



- Delicate migration planning

Years	Beam Species and Energies
2016	High statistics Au+Au d+Au energy scan
2017	High statistics Pol. p+p at 510 GeV
2018	⁹⁶ Zr+ ⁹⁶ Ru isobar run
2019-20	7.7-20 GeV Au+Au (BES-2)
2021	No Run ?
2022-23	200 GeV Au+Au with upgraded detectors Pol. p+p, p+Au at 200 GeV
2024---	Program TBD

Summary

- RCF performed remarkably well during Run 16
- About the right size of effort for current requirements
- Needs for replacement of old hardware, new tape generation & resources needs for 2017 and beyond
 - difficult with level of current budget
- Plan being developed for migrating facility to state of the art computing room in 2021

Backup

Recommendations from 2014 S&T review

1. BNL is encouraged to resolve the HVAC problem at RACF as soon as possible.
2. The RACF and the detector collaborations should analyze the processing capacity required to perform the necessary production runs (in units of HS06*years) and compare to the available capacity of RACF. If additional capacity is required, a plan to acquire the necessary capacity by 2018 should be developed. This plan should be submitted to DOE by January 1, 2015.

Response to the 2014 RHIC S&T Review Committee Recommendation

“The RACF and the detector collaborations should analyze the processing capacity required to perform the necessary production runs (in units of HS06*years) and compare to the available capacity of RACF. If additional capacity is required, a plan to acquire the necessary capacity by 2018 should be developed. This plan should be submitted to DOE by January 1, 2015.”

Prepared by Michael Ernst (RACF), Jerome Lauret (STAR) and Christopher Pinkenburg (PHENIX)

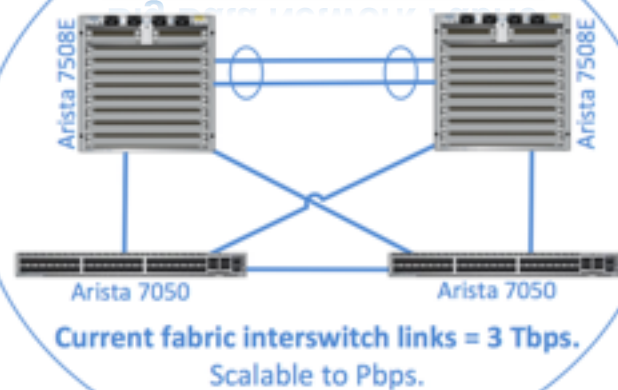
Table of Contents

Executive Summary	2
PHENIX Resource Requirements	7
PHENIX Data Analysis	7
Data Production and reconstructed output	8
sPHENIX	8
STAR Resource Requirements	11
Introduction	11
Input	11
Setting the basis for our requirements	12
Trend based calculation and adjustments	12
Data size projections	13
Event reconstruction projections	15
Straw-man plan	16
CPU capacity	16
CPU power adequacy – observations	18
Storage capacity	18
Summary	19
RACF Overview	21
Compute Farm	21
The HTCondor Local Resource Management System	22
Storage Infrastructure	22

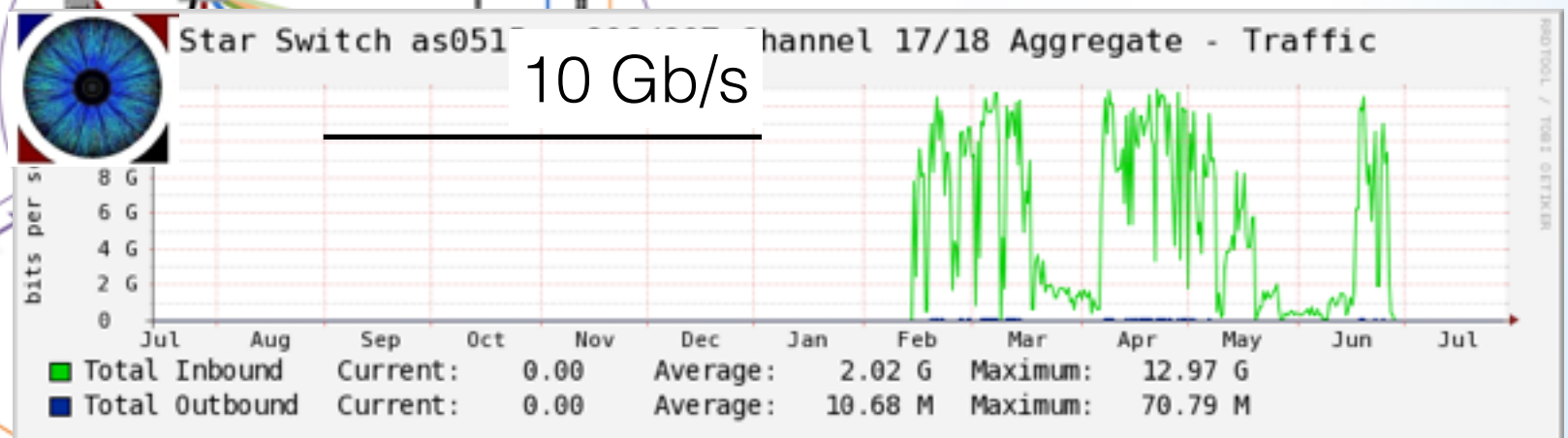
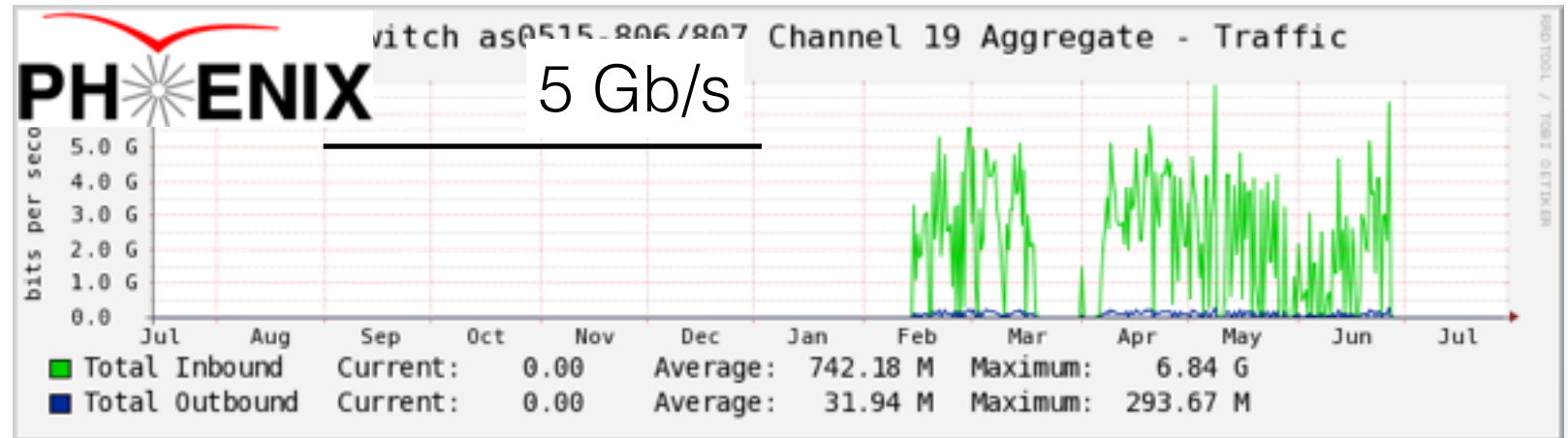
Performance in 2016

- No issue in data transfer from experiments to facility

Big Data Network Fabric

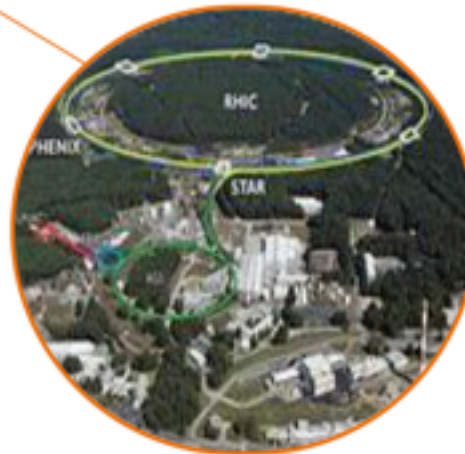


Tape Mass Storage



Disk storage = 25 PB

RHIC Facilities

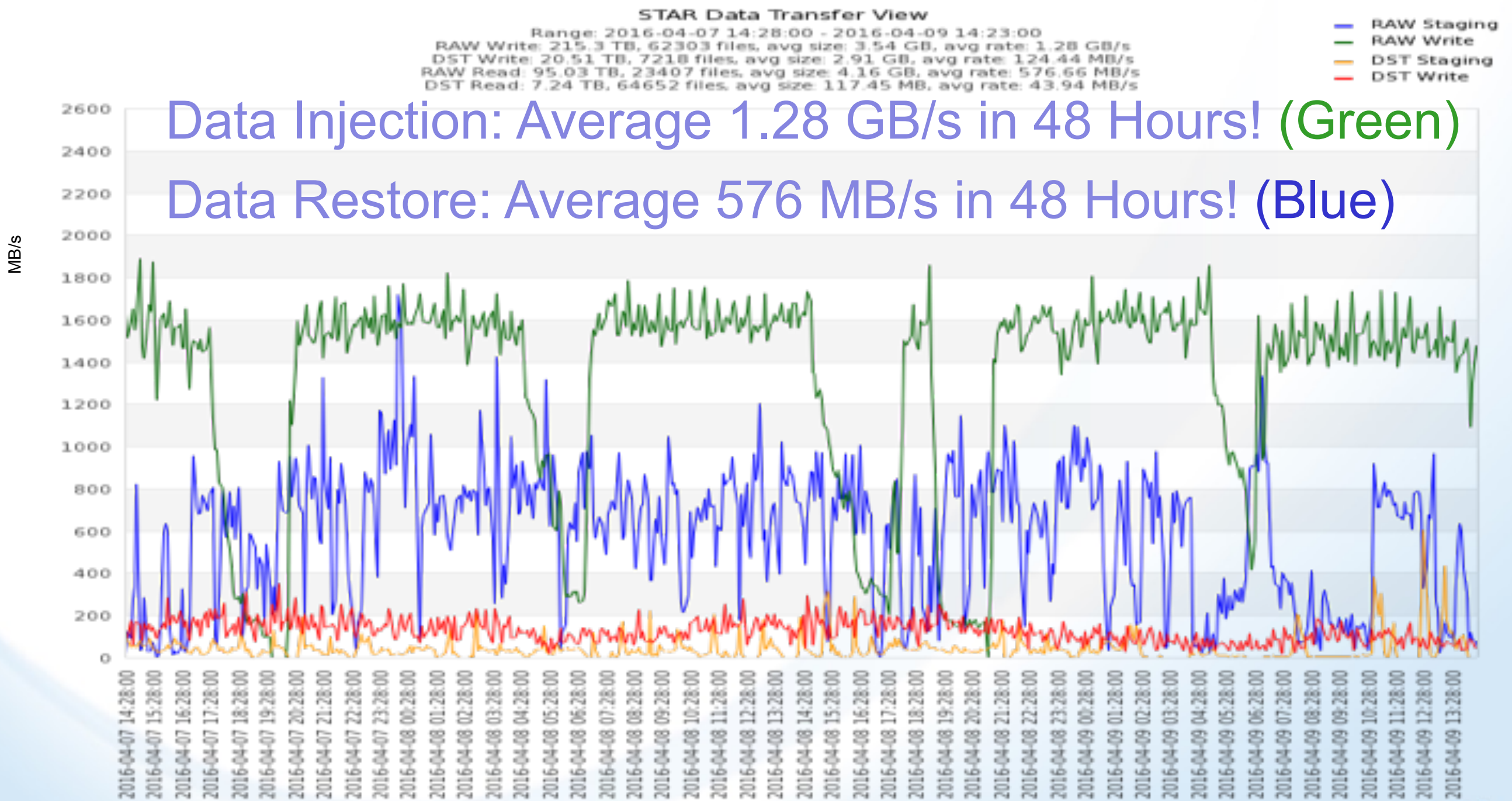


High Throughput Parallel Archiving

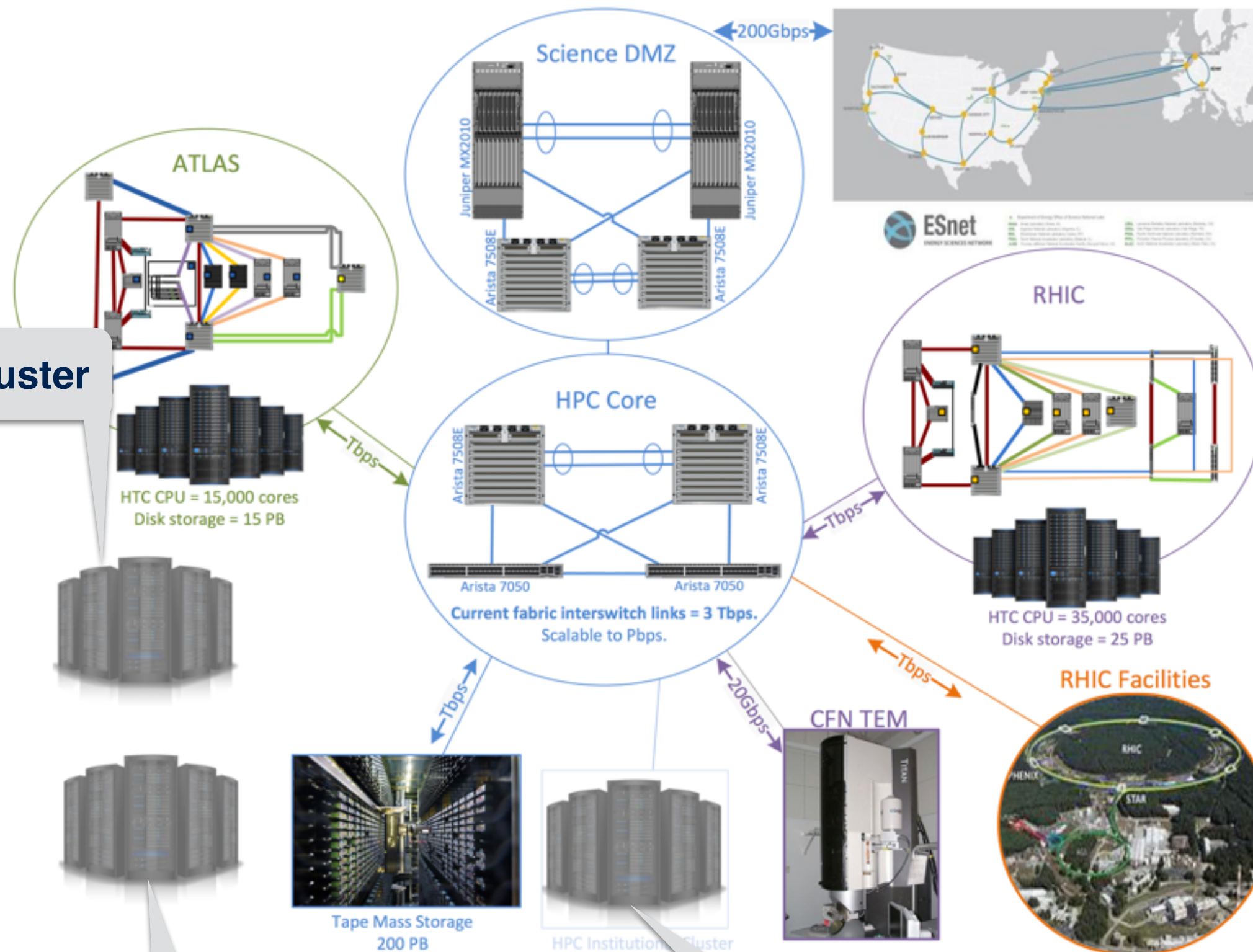
RHIC RUN 16 - STAR

Data Injection: Average 1.28 GB/s in 48 Hours! (Green)

Data Restore: Average 576 MB/s in 48 Hours! (Blue)



2017



CFR Design – An Incremental Approach

■ Power

- Day-one capability (2021) – 2.4 MW IT power (dedicated computing power). This is approximately double current RACF IT power.
- Provide provision for future 1.2 MW IT power increments to 6MW Max.

■ Cooling

- Day-one cooling capability to support 2.4 MW IT power
- Provide provision for future 1.2 MW IT power deployments

■ Space

- Day-one - Accommodate approximately 33% footprint expansion (Racks) within defined spaces.
- Day-one - Accommodate approximately 3,500 SF additional, unassigned space.
- Provide opportunity for future (long term) growth within the balance of the 725 facility. Both computing and offices.